# Machine Learning Operations Semester Project-Task 3

### Environmental Monitoring and Pollution Prediction System



**Student Name**: Moeed Asif
**Roll No**: 21i-0483
**Section** : Mlops-B

# 1. Introduction

This report outlines the implementation of a monitoring environment for a machine learning (ML) model served through a Flask API. The primary goal is to ensure real-time monitoring and observability of the model's performance and system metrics. Prometheus and Grafana were used to set up a robust monitoring stack to collect, query, and visualize metrics.

## 2. Objectives

The key objectives of this setup are:

- Monitor the ML model's inference performance, including latency and throughput.
- Track resource utilization (CPU, memory, etc.) of the Flask API hosting the model.
- Enable real-time visualization and alerts for critical metrics.
- Provide insights into system health and model efficiency.

## 3. Methodology

## 3.1. Task 3 Part 1-Set up Monitoring

The monitoring environment is built around the following components:

- **Flask API:** Serves the ML model and exposes metrics endpoints.
- **Prometheus**: Collects metrics from the Flask API and the system.
- **Grafana:** Visualizes the metrics and sets up alerts.

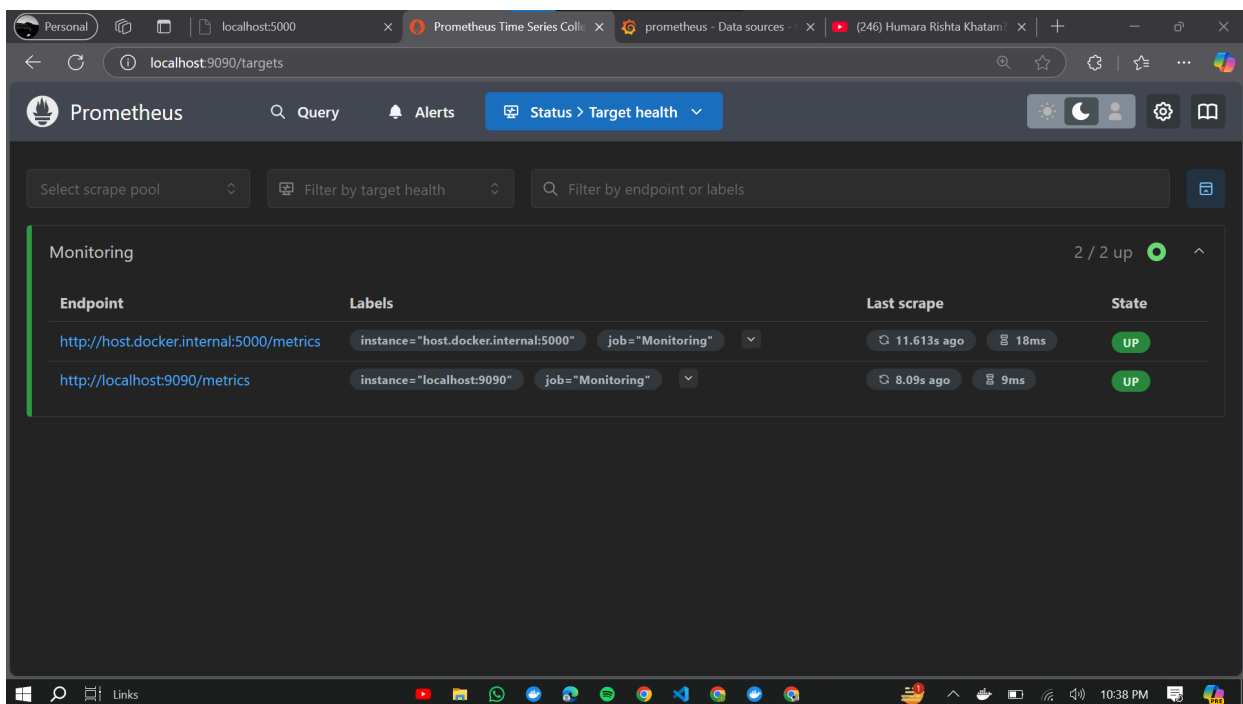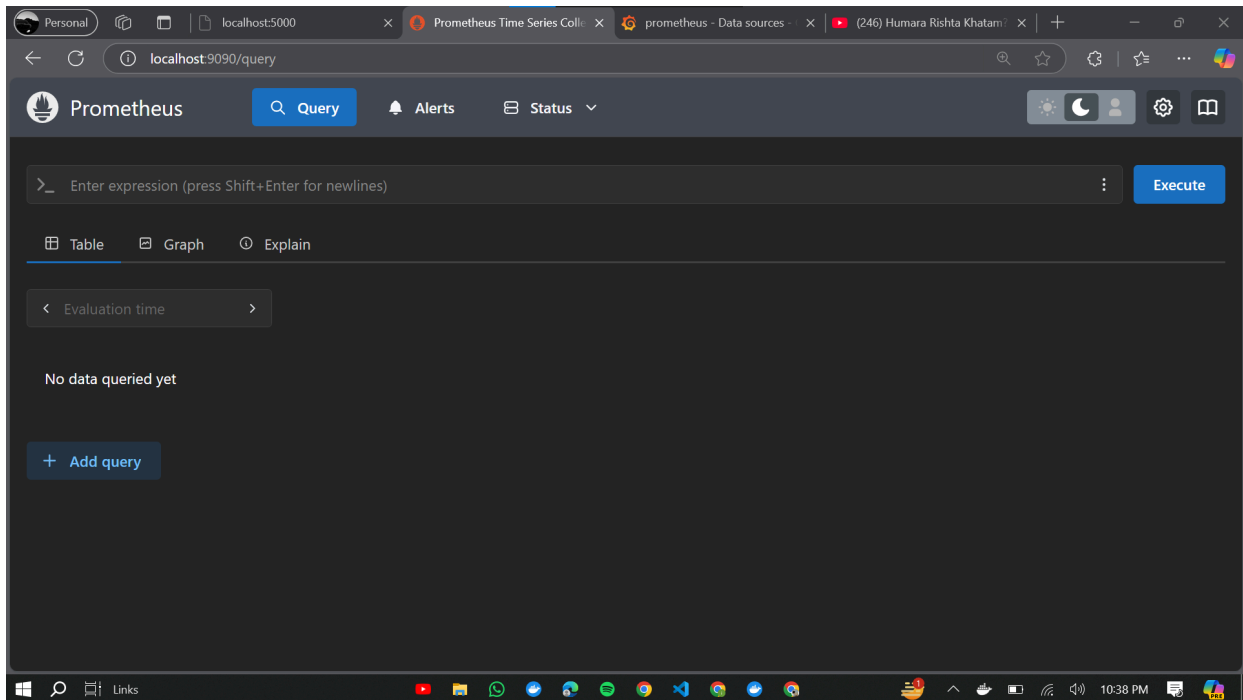## 3.2. Prometheus Setup

1. **Installing Prometheus:**
   - Prometheus was installed and configured on the server hosting the Flask API.
   - The prometheus.yml configuration file was updated to scrape metrics from the Flask API.
2. **Exposing Metrics:**
   - The Flask API was integrated with the prometheus-flask-exporter library to expose metrics.
   - Custom metrics were defined to monitor ML model-specific parameters, such as inference latency, request count, and error rates.
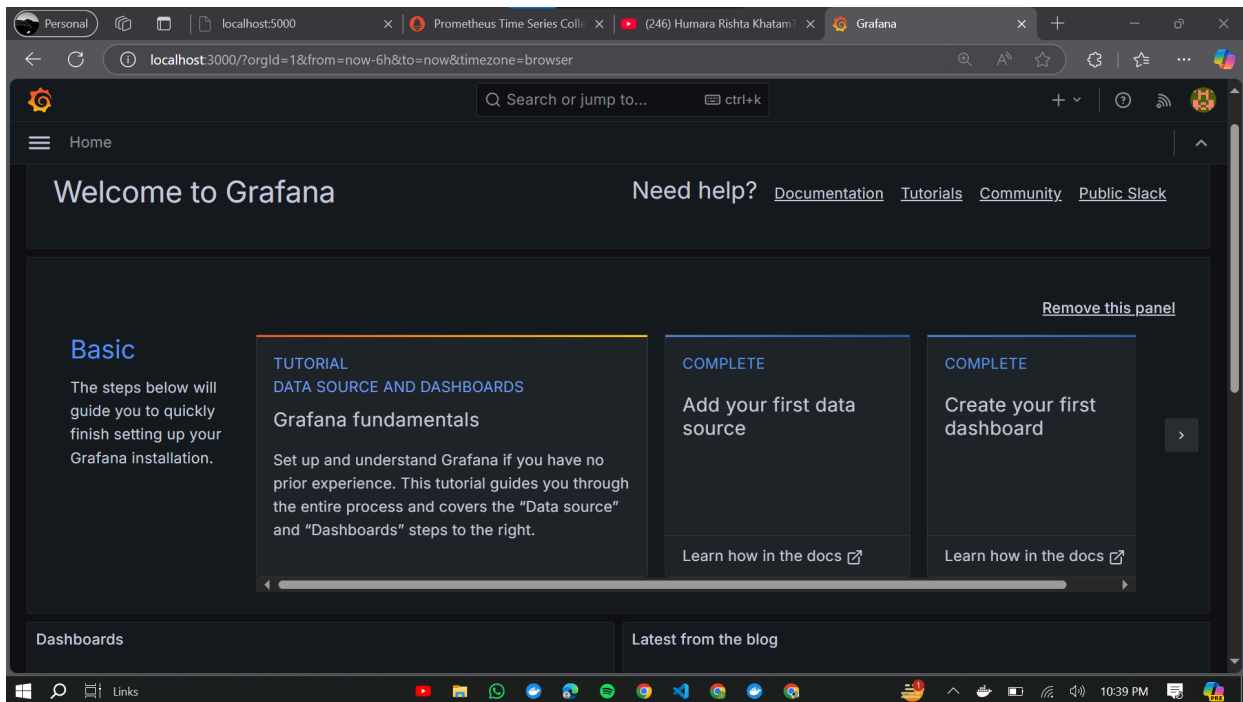3. **Scraping Metrics:**
   - Prometheus was configured to scrape metrics at regular intervals from the Flask API and system exporters (e.g., node_exporter for system metrics).
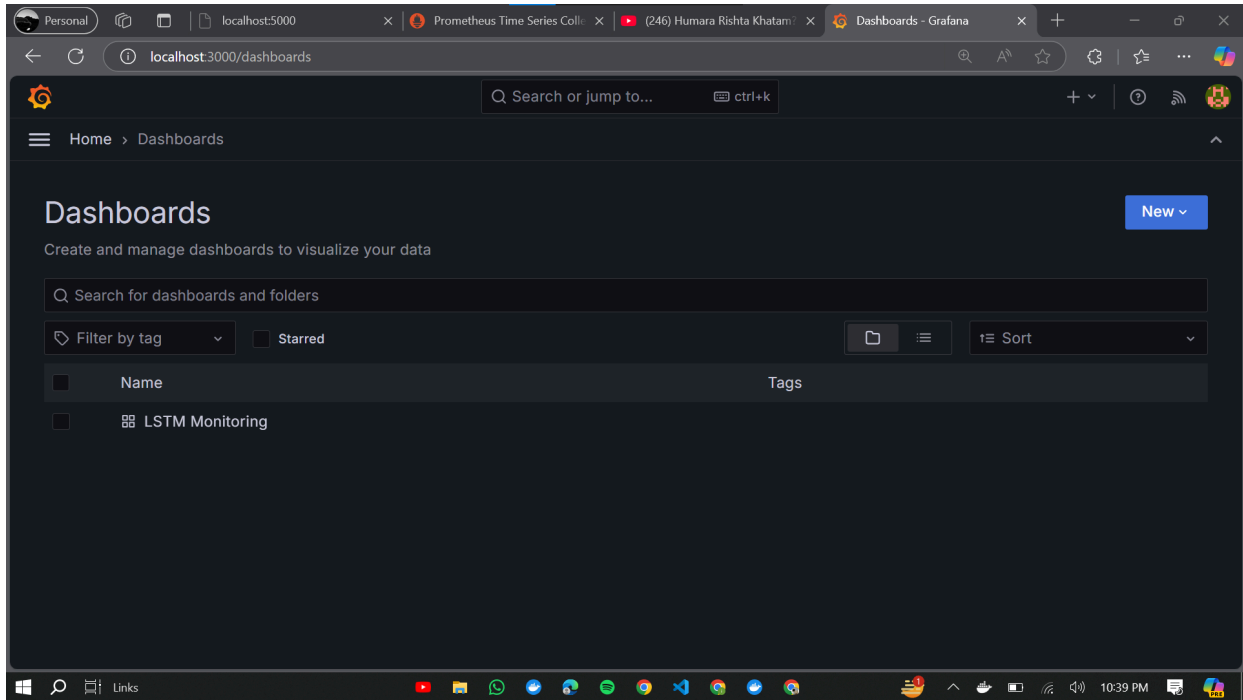
## 3.3. Grafana Setup

1. **Installing Grafana:**

- Grafana was installed on the same server as Prometheus for easy access.

2. **Connecting to Prometheus:**
   - A Prometheus data source was added in Grafana to query metrics.

3. **Creating Dashboards:**
   - Custom dashboards were created to display:
     - **API Metrics:** Request rates, latency, and error counts.
     - **System Metrics:** CPU usage, memory consumption, and disk IO.
     - **Model Metrics:** Prediction distribution, average inference time, and error rates.

## 4. Setting Up Alerts:
   - Alerts were configured for critical conditions, such as high latency or high error rates, to ensure timely responses.
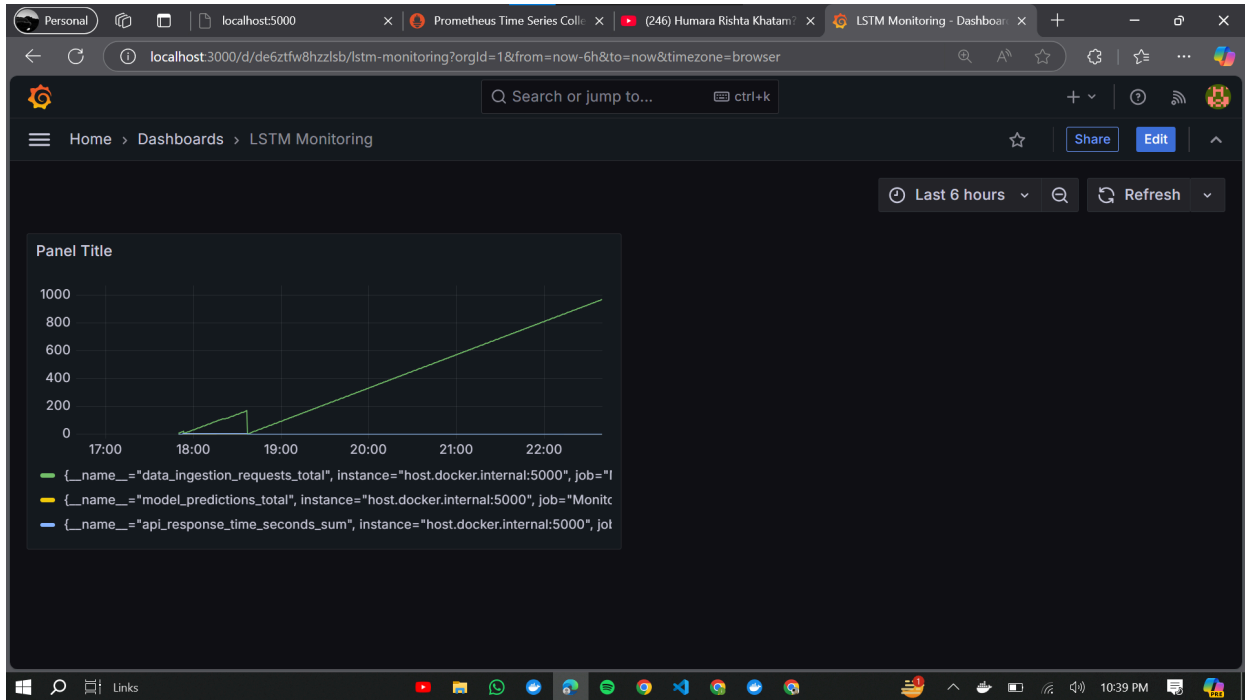
## 4. Metrics Monitored
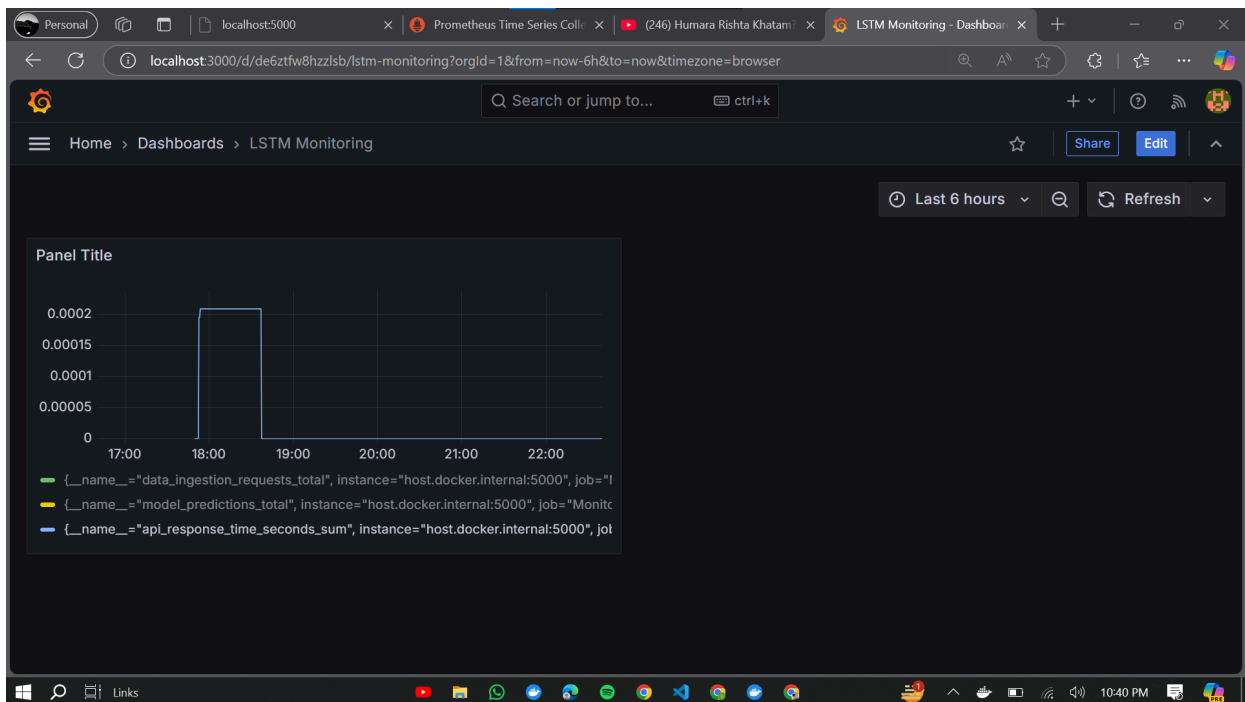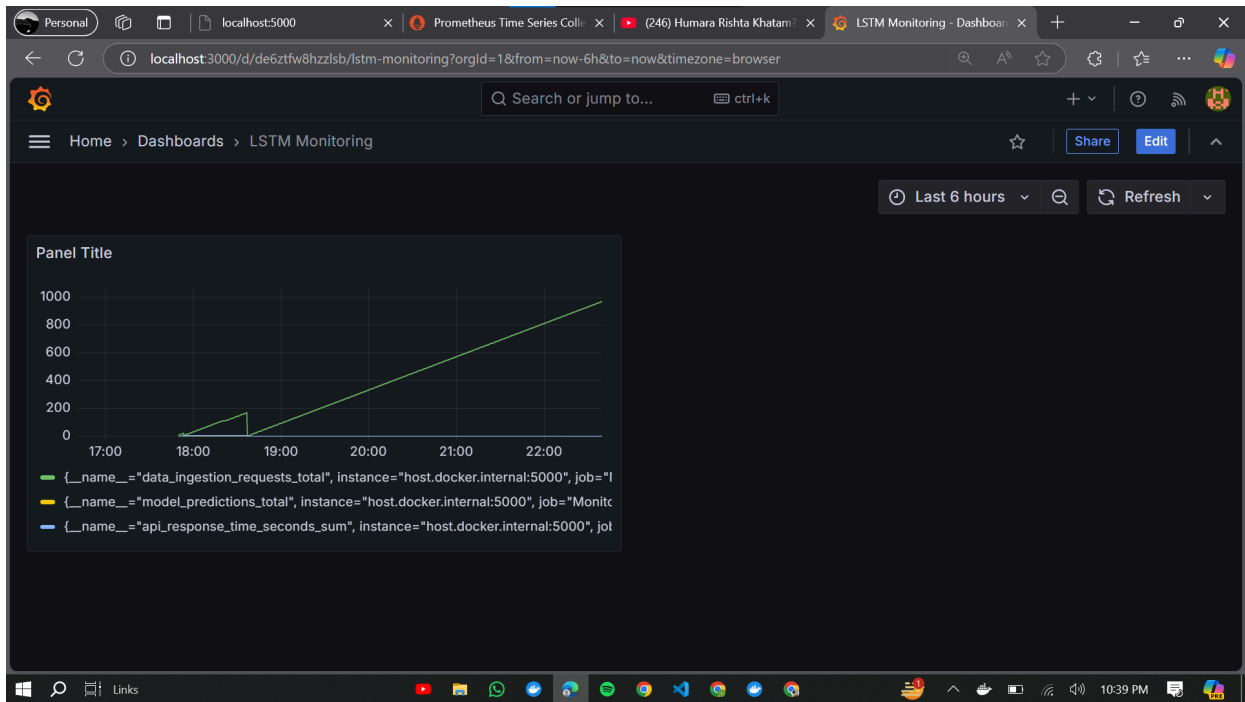
The following metrics were monitored:

- **API Metrics:**
  - Total requests served.
  - Request latency (p95, p99 percentiles).
  - Success and error rates.
- **System Metrics:**
  - CPU and memory usage.
  - Disk utilization.
  - Network IO.
- **Model-Specific Metrics:**
  - Inference time per request.
  - Distribution of predictions.
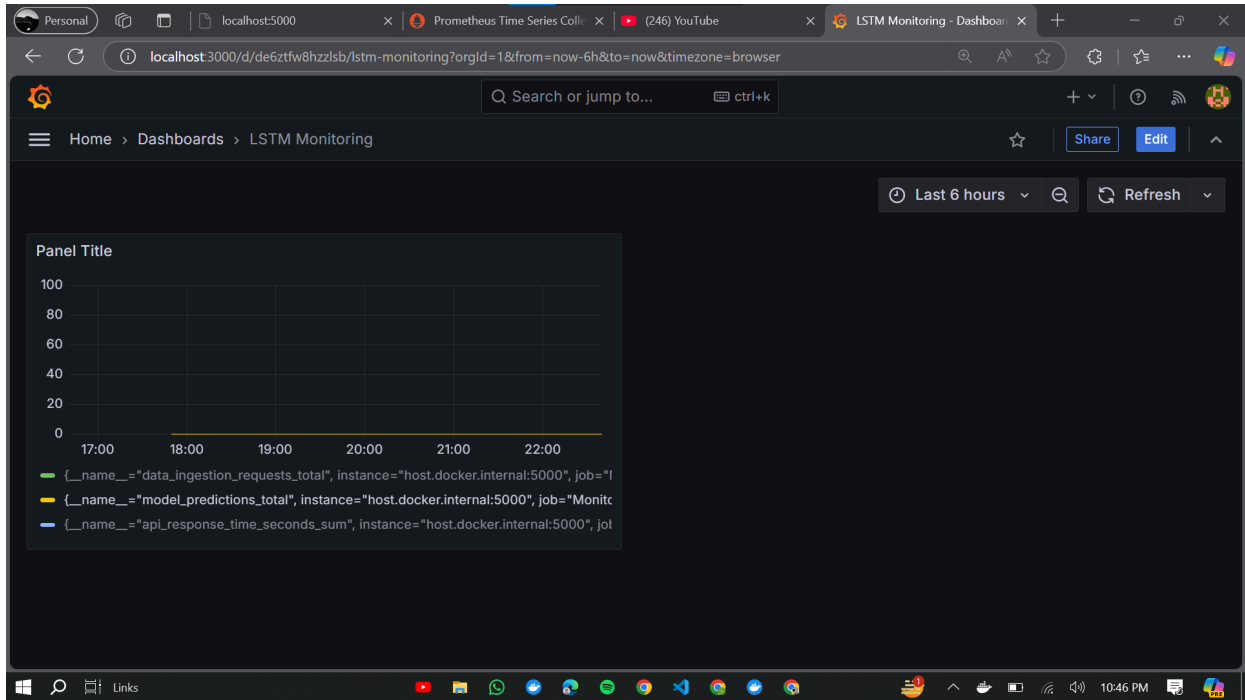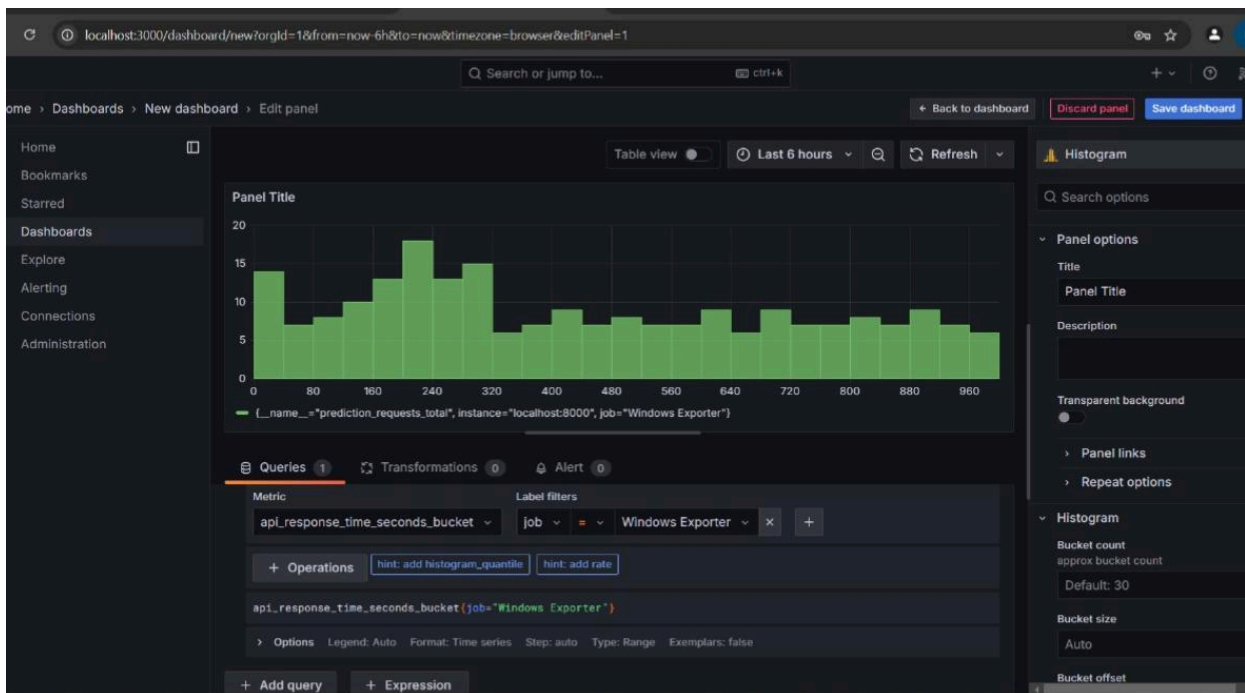  - Error rates based on model outputs.

## 5. Results

The monitoring setup successfully provided real-time visibility into the performance of the Flask API and the ML model. Key observations include:
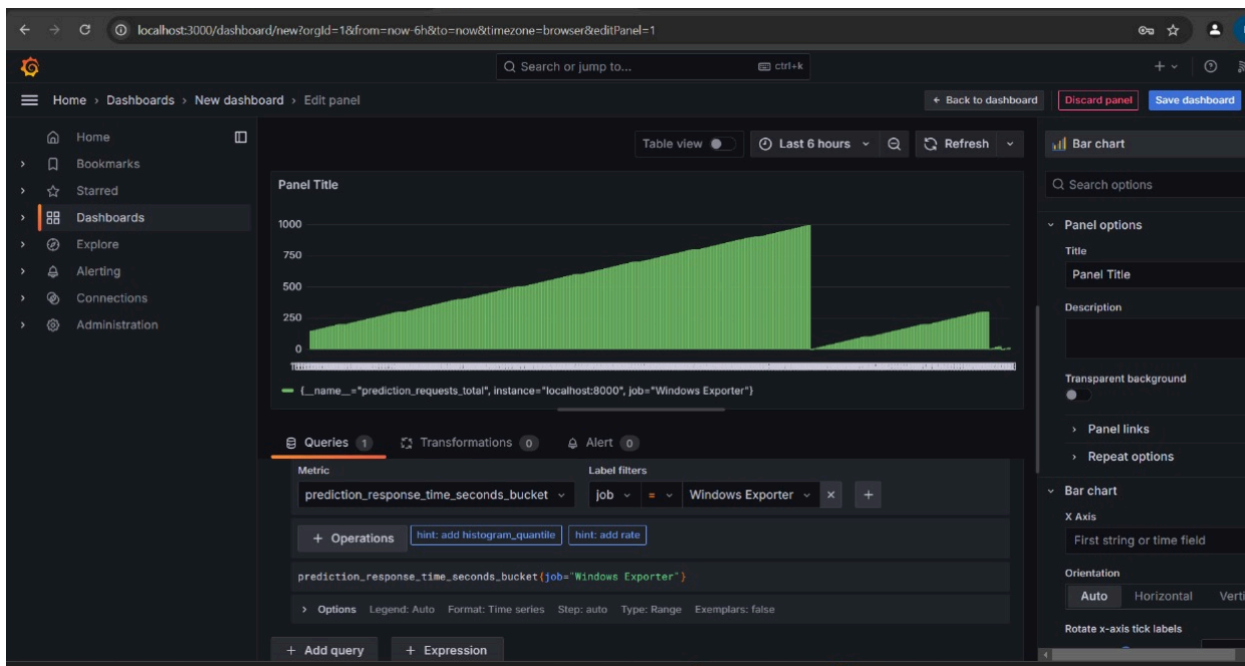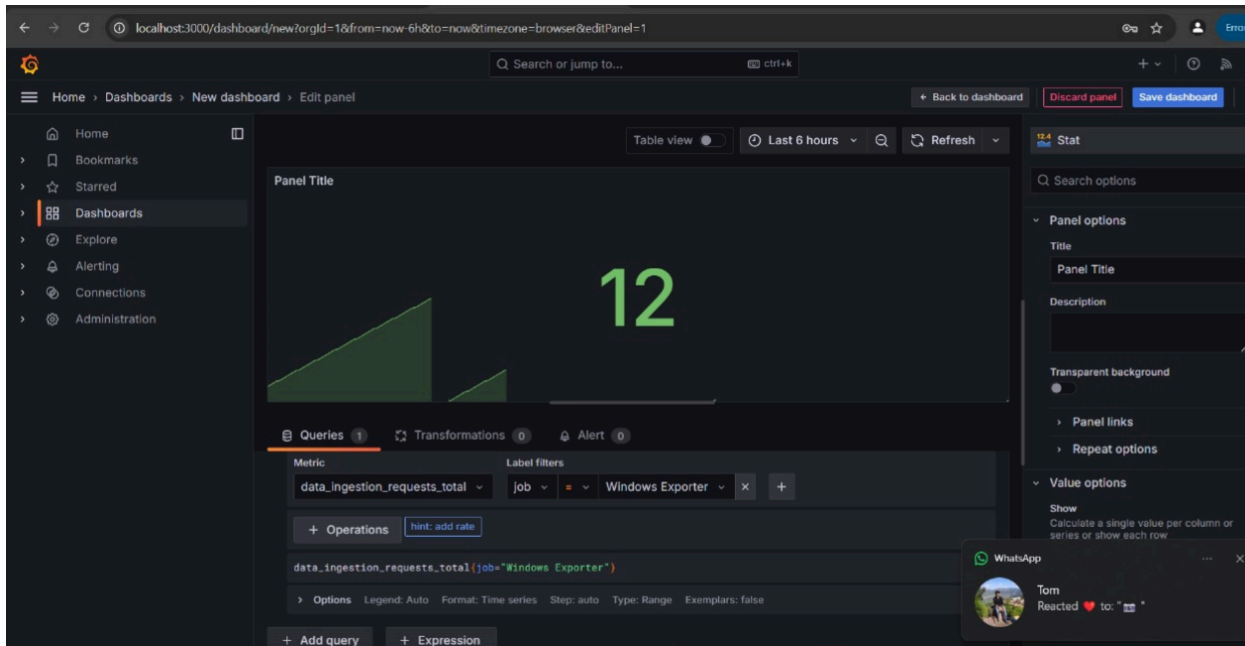
- Consistent response times under normal loads, with latency spikes during high traffic periods.
- Memory usage correlated with the size of input data processed by the model.
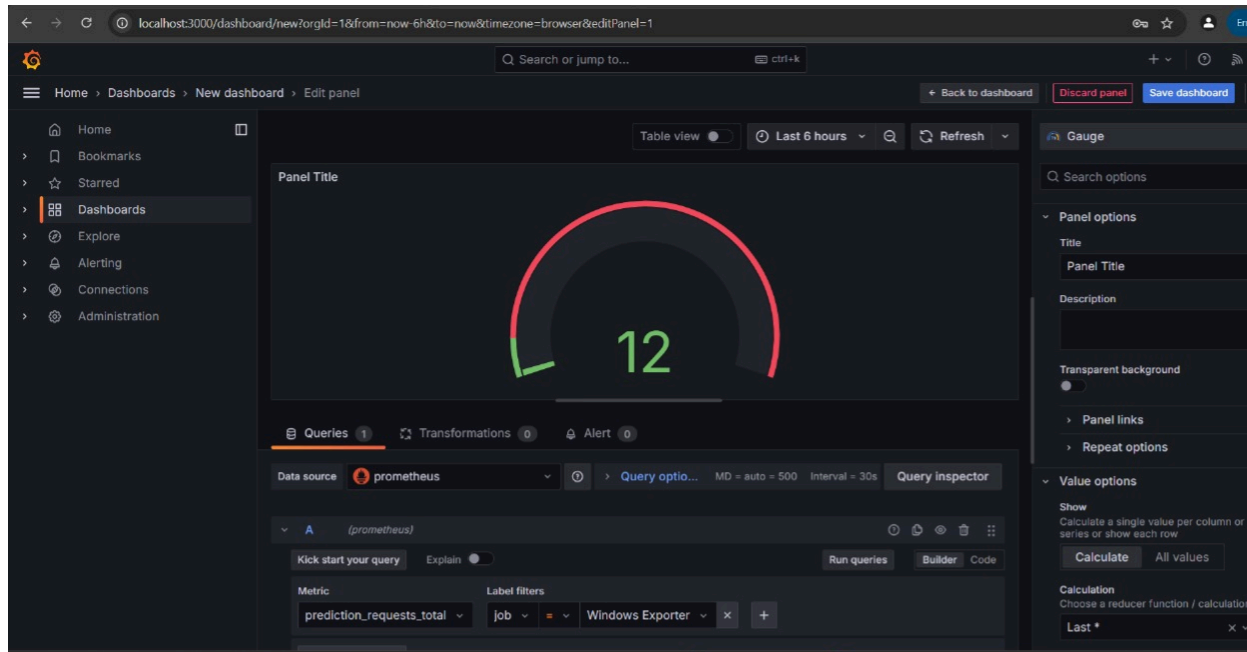- Alerts for latency and error thresholds effectively notified stakeholders of potential issues.
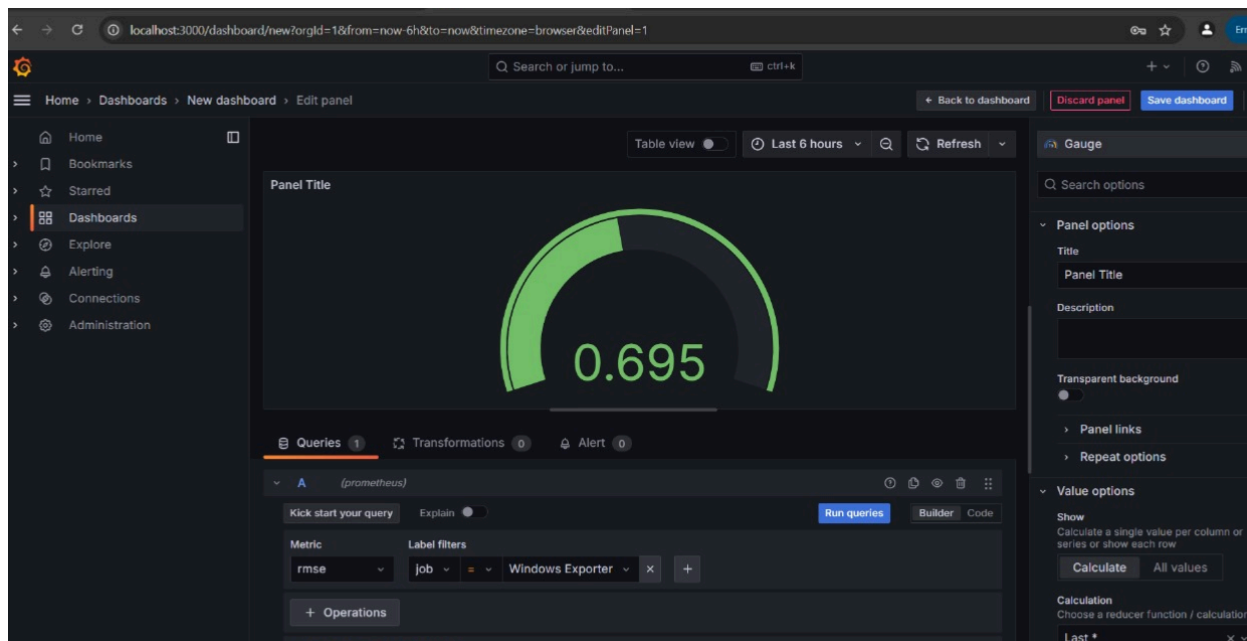
## 6. Task 3 Part 2-Performance

## 7. Task 3 Part 3 - Live Data Performance



## 8. Conclusion

The integration of Prometheus and Grafana provided an effective monitoring solution for the Flask API serving the ML model. This setup enables data-driven decision-making and ensures

the reliability of the deployed system. The use of real-time visualization and alerting mechanisms has improved the operational efficiency of the system.